



**University of
Zurich^{UZH}**

**Zurich Open Repository and
Archive**

University of Zurich
University Library
Strickhofstrasse 39
CH-8057 Zurich
www.zora.uzh.ch

Year: 2021

Ultralong Oxford Nanopore reads enable the development of a reference-grade perennial ryegrass genome assembly

Frei, Daniel ; Veekman, Elisabeth ; Grogg, Daniel ; Stoffel-Studer, Ingrid ; Morishima, Aki ; Shimizu-Inatsugi, Rie ; Yates, Steven ; Shimizu, Kentaro K ; Frey, Jürg E ; Studer, Bruno ; Copetti, Dario

Abstract: Despite the progress made in DNA sequencing over the last decade, reconstructing telomere-to-telomere genome assemblies of large and repeat-rich eukaryotic genomes is still difficult. More accurate basecalls or longer reads could address this issue, but no current sequencing platform can provide both simultaneously. Perennial ryegrass (*Lolium perenne* L.) is an example of an important species for which the lack of a reference genome assembly hindered a swift adoption of genomics-based methods into breeding programs. To fill this gap, we optimized the Oxford Nanopore Technologies sequencing protocol, obtaining sequencing reads with a N50 of 62 kb—an unprecedented value for a plant sample. The assembly of such reads produced a highly complete (2.3 of 2.7 Gb), correct (QV 45), and contiguous (contig N50 and N90 11.74 and 3.34 Mb, respectively) genome assembly. We show how read length was key in determining the assembly contiguity. Sequence annotation revealed the dominance of transposable elements and repeated sequences (81.6% of the assembly) and identified 38,868 protein coding genes. Almost 90% of the bases could be anchored to seven pseudomolecules, providing the first high quality haploid reference assembly for perennial ryegrass. This protocol will enable producing longer Oxford Nanopore Technology reads for more plant samples and ushering forage grasses into modern genomics-assisted breeding programs.

DOI: <https://doi.org/10.1093/gbe/evab159>

Posted at the Zurich Open Repository and Archive, University of Zurich

ZORA URL: <https://doi.org/10.5167/uzh-205024>

Journal Article

Accepted Version

Originally published at:

Frei, Daniel; Veekman, Elisabeth; Grogg, Daniel; Stoffel-Studer, Ingrid; Morishima, Aki; Shimizu-Inatsugi, Rie; Yates, Steven; Shimizu, Kentaro K; Frey, Jürg E; Studer, Bruno; Copetti, Dario (2021). Ultralong Oxford Nanopore reads enable the development of a reference-grade perennial ryegrass genome assembly. *Genome Biology and Evolution*, 13(8):evab159.

DOI: <https://doi.org/10.1093/gbe/evab159>

1

2

3

41 Title

5

6

72 Ultralong Oxford Nanopore reads enable the development of a reference-grade perennial ryegrass genome

8

9

103 assembly

11

12

13

144

15

16

175 Authors and affiliations

18

19

206 Daniel Frei¹, Elisabeth Veekman², Daniel Grogg³, Ingrid Stoffel-Studer³, Aki Morishima⁴, Rie Shimizu-

21

22

237 Inatsugi⁴, Steven Yates³, Kentaro K. Shimizu^{4,5}, Jürg E. Frey¹, Bruno Studer^{3*}, Dario Copetti^{3,4*}

24

25

26

278 ¹ Agroscope, Research Group Molecular Diagnostics, Genomics and Bioinformatics, Müller-

28

29

309 Thurgastrasse 29, 8820, Wädenswil, Switzerland

31

32

3310 ² DLF Seeds S/A, Højerupvej 31, 4660, Store Heddinge, Denmark

34

35

3611 ³ Molecular Plant Breeding, Institute of Agricultural Sciences, ETH Zurich, Universitätstrasse 2, 8092,

37

38

3912 Zurich, Switzerland

40

41

4213 ⁴ Department of Evolutionary Biology and Environmental Studies, University of Zurich,

43

44

4514 Winterthurerstrasse 190, 8057, Zurich, Switzerland

46

47

4815 ⁵ Kihara Institute for Biological Research, Yokohama City University, 641-12 Maioka, Totsuka-ward,

49

50

5116 Yokohama, 244-0813, Japan

52

53

54

55

56

57

58

59

60

17

18

***Authors for Correspondence**

Bruno Studer and Dario Copetti, Molecular Plant Breeding, Institute of Agricultural Sciences, ETH

Zurich, Universitätstrasse 2, 8092, Zurich, Switzerland, +41446320157, bruno.studer@usys.ethz.ch,

dario.copetti@usys.ethz.ch

23

24

Abstract

Despite the progress made in DNA sequencing over the last decade, reconstructing telomere-to-telomere

genome assemblies of large and repeat-rich eukaryotic genomes is still difficult. More accurate basecalls

or longer reads could address this issue, but no current sequencing platform can provide both

simultaneously. Perennial ryegrass (*Lolium perenne* L.) is an example of an important species for which

the lack of a reference genome assembly hindered a swift adoption of genomics-based methods into

breeding programs. To fill this gap, we optimized the Oxford Nanopore Technologies sequencing

protocol, obtaining sequencing reads with a N50 of 62 kb – an unprecedented value for a plant sample.

The assembly of such reads produced a highly complete (2.3 of 2.7 Gb), correct (QV 45), and contiguous

1

2

3

434 (contig N50 and N90 11.74 and 3.34 Mb, respectively) genome assembly. We show how read length was

5

6

735 key in determining the assembly contiguity. Sequence annotation revealed the dominance of transposable

8

9

1036 elements and repeated sequences (81.6% of the assembly) and identified 38,868 protein coding genes.

11

12

1337 Almost 90% of the bases could be anchored to seven pseudomolecules, providing the first high quality

14

15

1638 haploid reference assembly for perennial ryegrass. This protocol will enable producing longer Oxford

17

18

1939 Nanopore Technology reads for more plant samples and ushering forage grasses into modern genomics-

20

21

2240 assisted breeding programs.

23

24

25

2641

27

28

2942

30

31

32

3343

34

35

36

3744

38

39

4045

41

42

43

4446

45

46

47

4847 **Key words**

49

50

5148 *Lolium perenne*, forage grasses, perennial ryegrass, genomics, genome assembly, Oxford Nanopore

52

53

54

5549

56

57

58

59

60

Downloaded from <https://academic.oup.com/gbe/advance-article/doi/10.1093/gbe/evab159/6319026> by University of Zurich user on 14 July 2021

Significance statement

Sequencing eukaryotic genomes with long-read sequencing platforms is allowing to obtain genome assemblies of unprecedented quality also for many non-model organisms. However, especially in genomes with a high amount of long repeats, completeness and contiguity are limited by the quality (accuracy and/or length) of the input data. Here we present an innovative protocol for Oxford Nanopore Technologies' genomic plant DNA library preparation that considerably increases read length. We show how these exceptionally longer reads were key in obtaining a perennial ryegrass genome assembly with unprecedented statistics, both within its genus and among other plants of similar complexity. This work makes available a highly complete and contiguous genome assembly and the laboratory protocol necessary to produce long read data.

1

2

3

466

5

6

767

8

9

10

1168

12

13

14

1569

16

17

1870

19

20

21

2271Introduction

23

24

25

2672Over the last three years, the sequencing performance of the Oxford Nanopore Technologies (ONT)

27

28

2973platform has dramatically increased in terms of sequencing yield, accuracy, and read length ((Shafin et al.

30

31

32742020), <https://bit.ly/3peI7xH>, last accessed May 2021). Human or bacterial substrates can nowadays be

33

34

3575sequenced at read length N50 (the size of the shortest read that sums up to 50% of the total bases) that

36

37

3876reaches or goes beyond 100 kb (<https://bit.ly/2Ydjig4>, last accessed May 2021). As plant DNA is typically

39

40

4177more difficult to purify and preserve at high quality, plant sequencing performance metrics lag compared

42

43

4478to more accessible substrates. The current longest published plant ONT data (Lang et al. 2020; Zhou et al.

45

46

47792020) has read N50 values around 30 kb. Though this length surpasses most of the largest repeated

48

49

5080sequences (especially retroelements and centromeric repeats), they are still too short to resolve arrays of

51

52

5381such sequences.

54

55

56

57

58

59

60

5

<http://mc.manuscriptcentral.com/gbe>

Alone or in mixture with legumes, *Lolium* and *Festuca* spp. are the main crop used as a feed source for livestock. Perennial ryegrass (*Lolium perenne* L.) is the most cultivated grass species in Western European grasslands (Wilkins & Humphreys 2003). It is a diploid ($2n = 2x = 14$) species with a haploid genome of about 2.6 Gb (Kopecký et al. 2010). Like other forage grasses, perennial ryegrass is an outcrosser, meaning it retains high levels of heterozygosity. Heterozygosity and a high content in repetitive sequences are the main constraints that still hamper the development of a high-quality assembly for a forage grass (Byrne et al. 2015; Honig et al. 2016; Knorst et al. 2019). On the other hand, when striving to assemble heterozygous genomes, the two haplotypes assemble separately in allelic sequences, resulting in a diploid assembly (Copetti et al. 2021). Though it represents the true content of the nucleus of a heterozygous organism, as is, such assembly is unsuitable to be used as a haploid reference for variant calling. The closest species to the *Festuca* and *Lolium* species complex with a chromosome-scale assembly is orchardgrass (*Dactylis glomerata* L. (Huang et al. 2020)), but its incomplete nature (1.84 Gb out of ~2.6 Gb) and distant taxonomical placement (subtribe Dactylidinae) make it unsuitable to be used as a reference. The most proximal highly complete assemblies are *Brachypodium distachyon* and barley (*Hordeum vulgare* L.), with which perennial ryegrass shared the last common ancestor about 30 million years ago (Wu & Ge 2012). The availability of a haploid or homozygous perennial ryegrass genotype would simplify the assembly of a reference genome. Haploid or doubled haploid individuals can serve such task and represent invaluable material for the development of customized breeding populations.

1
2
3
4 100 To compensate for the absence of a high-quality reference assembly in forage grasses, we
5
6
7 101 sequenced Kyuss, a double haploid *L. perenne* genotype. The genotype derived from in vitro anther
8
9
10 102 culture (Begheyn et al. 2017). The application of an optimized DNA extraction (Russo et al. 2021) and the
11
12
13 103 development an improved sample handling protocol allowed to preserve DNA integrity, resulting in
14
15
16 104 unprecedented read lengths for a plant sample. *De novo* assembly of the sequence data resulted in a highly
17
18
19 105 complete, contiguous, and accurate genome assembly. This dataset can serve as a pivotal reference
20
21
22 106 assembly for genome-based studies in *Lolium* and *Festuca* biology and breeding.
23
24
25
26 107
27
28

29 108 **Results and Discussion**
30
31
32

33 109 To achieve high contiguity and completeness of the Kyuss genome assembly, we optimized the standard
34
35
36 110 ONT library preparation protocol. Here we report the most consequential modifications, while the
37
38
39 111 complete protocol used in this work is available as Supplementary protocol. The salient original
40
41
42 112 optimizations are the following: reducing DNA mechanical shock during library preparation; allowing for
43
44
45 113 longer elution times; and flushing the flow cell and reloading a second aliquot of library. We noticed that
46
47
48 114 mixing the components by tapping and not flicking the tube helped preserving the large molecules, likely
49
50
51 115 by avoiding DNA shearing and the accumulation of shorter fragments. Extending the elution time from
52
53
54 116 the AMPure beads resulted in a maximization of DNA recovery, particularly preserving the high
55
56
57
58
59
60

molecular weight fraction. Flushing the flow cell with Wash Solution countered pore clogging, restoring its productivity close to the initial levels and allowing to sequence more substrates on the same device. Also, we experienced that allowing for a 30 minutes incubation time before starting the run resulted in more active sequencing pores.

Upon base calling of the raw data, the baseline dataset for all the downstream analyses consisted of 69.6 Gb of ONT data in 2,061,375 reads, having a mean read length and N50 of 33.7 and 62.6 kb, respectively. The mean base call quality was QV 10.3. The prevalence of long reads was even more clear by considering other metrics: only 6.4% of the reads are longer than 100 kb but contain about 25% of the total bases. Inspecting the current literature, such metrics are unprecedented for a plant sample.

The size of the haploid Kyuss genome was estimated by flow cytometry and by counting k-mers from short-read data. The *in silico* estimation returned a value of 2.467 Gb, while flow cytometry estimated it at 2.720 Gb. For the sake of consistency with previous measurements in plants, we adopt the flow cytometry value as the estimated genome size of Kyuss. Furthermore, given that when compared to the tomato control, the flow cytometry peak profile was the same as its parent DH 6-47 (a diploid genotype (Begheyn et al. 2017), Figure 1a and b), we concluded that Kyuss is a doubled haploid plant.

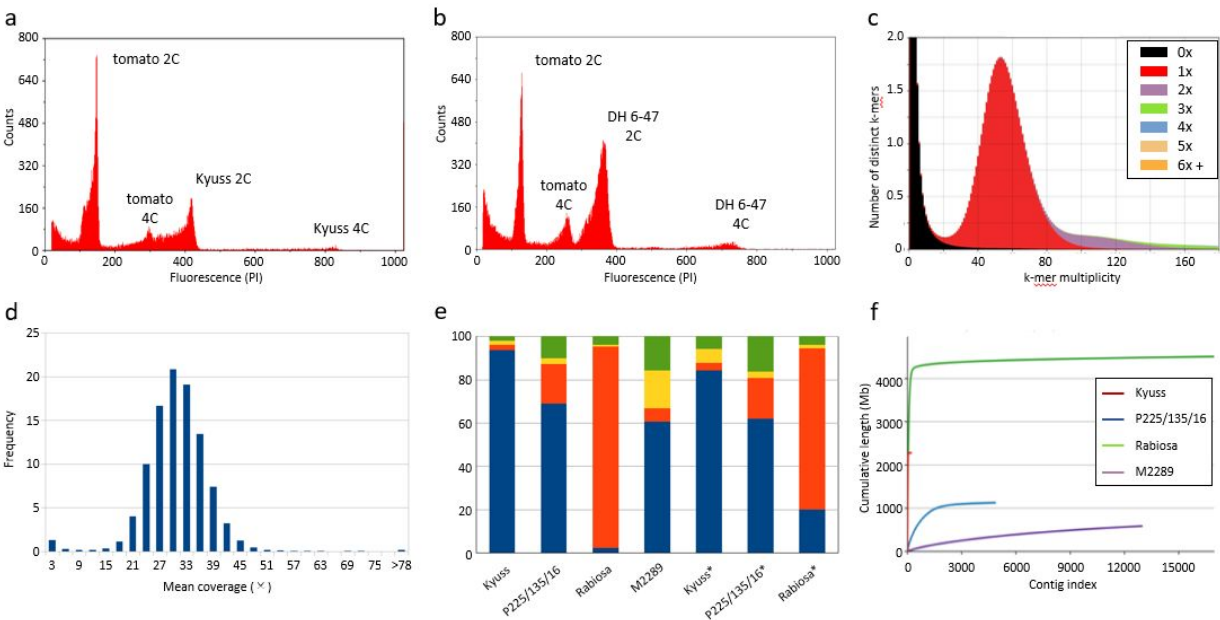


Fig. 1. Features of the Kyuss genome and assembly. a) Flow cytometry trace of Kyuss nuclei showing the

occurrence of peaks at the same position as in the diploid parent (b) when compared to the tomato external

standard. c) KAT comp spectra describing the occurrence of mostly single-copy k-mers (red area) in the

assembly under the main peak. The assembly is overall very complete (very small 0x area at multiplicities

higher than 20) and the repeated sequences are correctly represented (purple and green enrichment values

at multiples of the main peak). d) long read coverage distribution upon read alignment. The lack of

shoulders or additional peaks results entails the lack of collapsed or allelic regions in the assembly. e)

BUSCO analysis of the Kyuss and other public ryegrass assemblies ((Byrne et al. 2015), Copetti et al.

BioRxiv, (Knorst et al. 2019)). Kyuss assembly shows the highest completeness in terms of conserved

single-copy orthologs (SCOs), with only 4% of the models being fragmented or missing. The ‘Rabiosa’

assembly is a diploid assembly, thus most of the SCOs are expected to be identified twice. The columns

with the asterisk denote BUSCO scores for the predicted gene models. Blue: single copy, orange: duplicated, yellow: fragmented, green: missing models. f) Total size and contiguity of the ryegrass assemblies evaluated by cumulative sequence length. The expected total size of the assemblies is around 2,500-2,700 Mb, except for Rabiosa where the diploid assembly should result in approximately 5,200 Mb. The high contiguity of the Kyuss assembly is denoted by the sharp vertical raise of the contig index approaching rapidly the total assembly size. In comparison, the ‘P226/135/16’ and ‘M2289’ assemblies show a dramatically lower completeness and contiguity.

The *de novo* assembly of the long-read data resulted in 1,920 contigs spanning 2.28 Gb, with N50 and N90 values of 11.74 and 3.34 Mb (the 60 and 201 longest sequences), respectively. Upon correction of one chimeric contig and short-read polishing, the continuity metrics decreased slightly to 11.28 and 3.32 Mb in 1,935 sequences (Table 1). The base-level accuracy was estimated to be QV 45. By intersecting k-mer abundance in the raw reads and in the assembly, the completeness was estimated at 99.39% (Figure 1c). Upon realignment of the long reads to the assembly, the unimodal coverage distribution entailed absence of collapsed genomic regions in the assembly (Figure 1d). The alignment rate was 99.89% and 89.89% for the long and short reads, respectively. Similarly, long and short reads mapped uniquely to 99.71% and 89.36% of the assembly bases. The overall lower mapping rate and breadth of

1
2
3
4 161 coverage observed with the short read data can be explained by ambiguous alignments due to repeats and
5
6
7 162 to a sequencing bias against extreme GC content regions seen with the Illumina platform (Browne et al.
8
9
10 163 2020). The assembly spanned 83.9% of the estimated haploid genome size and contained the vast majority
11
12
13 164 of the single-copy orthologs expected to be present (93.5% in single copy, 2.5% in more than one copy),
14
15
16 165 leaving only 2% of the models fragmented and 2% missing (Figure 1e). When compared to the publicly
17
18
19 166 available ryegrass assemblies, Kyuss showed much higher completeness and contiguity than *L. perenne*
20
21
22 167 ‘P226/135/16’ (Byrne et al. 2015) and the *L. multiflorum* ‘M2289’ (Knorst et al. 2019) assemblies (Figure
23
24
25 168 1f, Table 1). While the Italian ryegrass ‘Rabiosa’ assembly also reached scaffold N50 values in the multi-
26
27
28 169 Mb range (N50 2.9 Mb, (Copetti et al. 2021)), N50 is still ~4 times shorter than Kyuss. Furthermore,
29
30
31 170 Rabiosa’s diploid nature makes it not suitable to be used as a haploid reference for mapping-based
32
33
34 171 analyses (e.g. SNP and SV calling, phasing). In conclusion, the Kyuss assembly is a highly complete,
35
36
37 172 accurate, and contiguous representation of a haploid genome.
38
39
40
41 173 To estimate the impact of read length on assembly completeness and contiguity, we created sets of shorter
42
43
44 174 input reads by cutting the sequence at different positions, thus maintaining the same number of input
45
46
47 175 bases. Though with different input read lengths the total assembly size did not change, when the reads
48
49
50 176 were shorter, all contiguity values decreased (Supplementary Figure 1, Supplementary Table 1). Notably,
51
52
53 177 the effect was more pronounced when introducing a second cut (C20+20 dataset, contig N50 1.9 Mb and
54
55
56 178 N90 0.5 Mb) rather than when the cut was single and more internal in the read (C30, C40, C60 datasets,
57
58
59
60

contig N50 of 6.6 Mb and N90 1.8 Mb). This can be consequence of the fact that in the C20+20 set more reads were cut to a shorter size, while in the other datasets less reads were long enough to be cut further inwards. These findings confirm the primary importance of read length in obtaining larger contigs by providing connectivity between adjacent genomic regions.

	Kyuss*	P226/135/16	Rabiosa*	M2289
	<i>L. perenne</i>	<i>L. perenne</i>	<i>L. multiflorum</i>	<i>L. multiflorum</i>
Reference	This study	(Byrne et al. 2015)	(Copetti et al. 2021)	(Knorst et al. 2019)
Est. genome size (Gb)	2.720	2.068	2.464	2.500
Assembly size (Gb)	2.281	1.128	4.531	0.585
% of genome assembled	83.9	54.6	183.9	23.4
# of sequences	1,935	48,415	226,949	129,579
N50 (kb)	11,276	70	2,941	5

N90 (kb)	3,320	14	283	2
L50 (#)	65	4,908	443	37,162
L90 (#)	209	16,951	1,984	103,446

*statistics of the contigs/scaffolds before being placed on pseudomolecules

Table 1. Statistics of the *L. perenne* Kyuss genome assembly and comparison with other public ryegrass

assemblies. The fraction of the assembled genome is based upon the genome size estimation provided in the respective studies. In the Rabiosa assembly, most of the allelic regions are represented as separate sequences, thus reaching the diploid genome size.

Though sequencing large and complex genomes is becoming easier and less expensive, optimizing the experimental conditions can help an optimal resource allocation. To estimate the coverage sufficient to give the best assembly possible given the current read length, assembly algorithm, and given the complexity and composition of the Kyuss genome, we assembled subsets of reads and inferred N50 values at higher coverages. Though total assembly size was not affected by the different coverage levels (Supplementary Table 2), with less input data (with same average read length), contig Nx values decreased rapidly (Supplementary Figure 2 a). We estimated that, at 52x and 70x coverage the contig N50 value has already reached 80% and 90%, respectively, of the maximum value (22.7 Mb) it could

theoretically reach (Supplementary Figure 2 b). These deducted N50 values allowed to estimate a sequencing depth threshold where the increase in contiguity is minimal – making sequencing additional bases less cost effective.

By exploiting genetic linkage of a *L. perenne* segregating population and collinearity with barley, we anchored 2 Gb (or 88% of the bases) of the assembly to seven chromosome pseudomolecules, corresponding to the haploid perennial ryegrass karyotype. Of the 235 contigs that constituted the pseudomolecules, 219 (85.5% of the bases) were oriented. Only 1,700 sequences spanning in total 274 Mb of sequence could not be anchored (N50 3.1 Mb, N90 83 kb). Besides such contigs being shorter, the unassignment to a pseudomolecule can be explained by the lack of Kyuss-specific markers (the segregating population was developed from parents that are not related to Kyuss) and by the occurrence of ryegrass-specific sequences with no counterpart in barley.

By similarity search, we identified that 60.0% of the assembly bases were composed of transposable elements and that 21.6% were constituted of uncharacterized repeated sequences. The gene annotation produced 38,868 protein-coding gene models, whose coding regions spanned 2% of the assembly. Completeness of the gene annotation set was estimated with BUSCO at 87,6% (Figure 1 e, (Seppey et al. 2019)) and 92,19% with the PLAZA core gene families (411 out of 7076 gene families

missing). Both values are on par with another perennial ryegrass gene annotation (Van Bel et al. 2012; Blanco-Pastor et al.).

Here we present a method that yields exceptionally long ONT reads for a plant sample and show how read length is key in obtaining highly contiguous and complete assemblies of repeat-rich plant genomes. The resulting assembly is the first reference-grade genome map for perennial ryegrass. Its availability will enable genomics-assisted development of breeding programs.

Materials and methods

For complete details, see supplementary materials and methods. Anther culture of the *L. perenne* genotype DH 6-47 (Begheyn et al. 2017) was used to establish double haploid individuals. One genotype (Kyuss) showing only one allele at three loci was selected for sequencing and clonally propagated. Ploidy level and genome size were measured via flow cytometry using tomato (*Solanum lycopersicum*) as a standard. The genome size was also estimated *in silico*, using short read data and Jellyfish (v2.4.2 (Marçais & Kingsford 2011)).

High molecular weight DNA was extracted from leaves according to (Russo et al. 2021) and the ONT sequencing library was prepared with an in-house optimized Genomic DNA by Ligation (SQK-LSK109, version GDE_9063_v109_revU_14Aug2019) protocol (Supplementary protocol). The library was

loaded into FLO-PRO002 flow cells and sequenced on a PromethION instrument (Oxford Nanopore Technologies Oxford OX4 4DQ, United Kingdom). About 200 ng genomic DNA were sheared to a mean fragment size of 500 bp and an Illumina TruSeq library was produced. The library was sequenced on a Novaseq 6000 (Illumina Inc, California, USA) instrument in 2×150 bp mode.

The ONT raw data was basecalled with Guppy (v4.0.14, <https://community.nanoporetech.com>, last accessed May 2021), keeping reads having a minimum q-score of 7. After adapter removal, sequences shorter than 2 kb were removed. The genome was assembled with Flye (v2.7.1-b1590 (Kolmogorov et al. 2019)). An auxiliary assembly was generated with Shasta (v.0.6.0 (Shafin et al. 2020)). The Flye assembly contigs were at first polished with the built-in Flye polishing module (twice), then twice with medaka (v1.1.1, <https://nanoporetech.github.io/medaka/>, last accessed May 2021), and lastly by two rounds of Pilon (v1.23 (Walker et al. 2014)). The haploid status of the assembly was confirmed by aligning long and short reads and plotting the coverage distribution. Sequence accuracy was measured from the short-read alignment file, parsed with the SAMtools stats subcommand. The assembly completeness was assessed with KAT (v2.4.2 (Mapleson et al. 2017)). The completeness of gene space was assessed with BUSCO (v3.0 (Seppey et al. 2019)). Contigs belonging to organelles were identified by aligning the assembly to *Lolium* chloroplast and mitochondrial deposited genomes. Such contigs were kept in the assembly and the sequences were flagged. Endosymbiont sequences were identified by BLASTN (BLAST+ suite, v2.9.0+ (Camacho et al. 2009)) against a collection of *Epichloe* and *Neotyphodium* sequences retrieved from NCBI nr.

1
2
3 252 Chromosome pseudomolecules were generated with ALLMAPS (v0.7.7 (Tang et al. 2015)) with a
4
5
6 253 genetic linkage map (Pfeifer et al. 2013) exploiting the collinearity with barley as input evidence. The *L.*
7
8
9 254 *perenne* ESTs associated with the genetic markers were aligned to the Kyuss assembly with TBLASTX.
10
11
12 255 Collinearity with barley (Mascher et al. 2017) was established aligning a draft Kyuss gene annotation
13
14
15 256 obtained by projecting *L. multiflorum* ‘Rabiosa’ gene models (Copetti et al. 2021) to the Kyuss contigs.
16
17
18 257 Upon aligning the two proteomes, the collinear blocks were determined with MCXScanX (v2 (Wang et al.
19
20
21 258 2012)). Chimeric contigs (containing stretches of more than four collinear genes aligning to a second barley
22
23
24 259 pseudomolecule) were split in two or more sequences.

25
26
27
28 260 Repeats and transposable element were identified with RepeatMasker (v4.0.6 (Smit A.F.A. et al.))
29
30
31 261 and a custom-built *Lolium* repeat library (Copetti et al. 2021). Protein-coding genes were annotated by
32
33
34 262 combining *ab initio* and homology-based evidence. The latter set was constituted of proteomes of
35
36
37 263 annotations of closely-related species (*Brachypodium*, barley, bread wheat (*Triticum aestivum* L.),
38
39
40 264 perennial and Italian ryegrass) and transcripts reconstructed from publicly-available RNA-Seq data from
41
42
43 265 NCBI SRA. Gene predictions with homology to transposable elements were removed prior to renaming the
44
45
46 266 models.

47
48
49
50 267
51
52
53
54 268 **Data Availability**

269 The raw sequencing reads are available at NCBI BioProject PRJNA690687, the assembly is available at
270 accession number JAEVFB000000000. The gene and repeat annotations are available at DOI:
271 10.25739/frsm-e984

272

273

274 **Supplementary Files**

275 Supplemental data (contains Supplementary materials and methods, supplemental protocol, supplemental
276 references)

277 Supplemental tables and figures

278

279 **Acknowledgements**

280 The authors are grateful to IT Support Group D-HEST, ETH Zurich for computational support, Roberto
281 Copetti for modeling support.

282 **Funding**

283 This work was funded by an ETH Zurich starting grant assigned to BS; by JST CREST [JPMJCR16O3],
284 Swiss National Science Foundation [31003A_182318], and URPP Evolution in Action to KKS.

1

2

3

4285

5

6

7286

8

9

10287

11

12

13288

14

15

16289

17

18

19290

20

21

22291

23

24

25

26292

27

28

29293

30

31

32

33294

34

35

36

37295

38

39

40

41296

42

43

44297

45

46

47

48298

49

50

51

52299

53

54

55

56

57

58

59

60

Author contributions

SY, BS, and DC conceived the experiment; DF developed the optimized ONT library protocol and produced the sequencing data; IS-S and AM performed the other molecular biology procedures; DG established the double haploid plants; EV and DC performed bioinformatic analyses; RS-I, KKS, JF, and BS supervised the analyses and provided critical feedback; DC wrote the manuscript. All authors read and approved the submitted version of the manuscript.

Literature cited

- 300 Begheyn RF, Roulund N, Vangsgaard K, Kopecký D, Studer B. 2017. Inheritance patterns of the response
301 to in vitro doubled haploid induction in perennial ryegrass (*Lolium perenne* L.). *Plant Cell Tiss Organ*
302 *Cult.* 130:667–679. doi: 10.1007/s11240-017-1255-y.
- 303 Blanco-Pastor JL et al. Canonical correlations reveal adaptive loci and phenotypic responses to climate in
304 perennial ryegrass. *Molecular Ecology Resources*. n/a. doi: <https://doi.org/10.1111/1755-0998.13289>.
- 305 Browne PD et al. 2020. GC bias affects genomic and metagenomic reconstructions, underrepresenting
306 GC-poor organisms. *GigaScience*. 9. doi: 10.1093/gigascience/giaa008.
- 307 Byrne SL et al. 2015. A synteny-based draft genome sequence of the forage grass *Lolium perenne*. *Plant J.*
308 84:816–826. doi: 10.1111/tpj.13037.
- 309 Camacho C et al. 2009. BLAST+: architecture and applications. *BMC Bioinformatics*. 10:421. doi:
310 10.1186/1471-2105-10-421.
- 311 Copetti D et al. 2021. Evidence for high intergenic sequence variation in heterozygous Italian ryegrass
312 (*Lolium multiflorum* Lam.) genome revealed by a high-quality draft diploid genome assembly. *bioRxiv*.
313 2021.05.05.442707. doi: 10.1101/2021.05.05.442707.
- 314 Honig J et al. 2016. ASM173568v1 - Genome - Assembly - NCBI.
315 https://www.ncbi.nlm.nih.gov/assembly/GCA_001735685.1/ (Accessed January 26, 2021).
- 316 Huang L et al. 2020. Genome assembly provides insights into the genome evolution and flowering
317 regulation of orchardgrass. *Plant Biotechnology Journal*. 18:373–388. doi:
318 <https://doi.org/10.1111/pbi.13205>.
- 319 Knorst V et al. 2019. First assembly of the gene-space of *Lolium multiflorum* and comparison to other
320 Poaceae genomes. *Grassland Science*. 0. doi: 10.1111/grs.12225.
- 321 Kolmogorov M, Yuan J, Lin Y, Pevzner PA. 2019. Assembly of long, error-prone reads using repeat
322 graphs. *Nature Biotechnology*. 37:540–546. doi: 10.1038/s41587-019-0072-8.
- 323 Kopecký D et al. 2010. Physical distribution of homoeologous recombination in individual chromosomes
324 of *Festuca pratensis* in *Lolium multiflorum*. *Cytogenet. Genome Res.* 129:162–172. doi:
325 10.1159/000313379.

1
2
3 326 Lang D et al. 2020. Comparison of the two up-to-date sequencing technologies for genome assembly: HiFi
4 327 reads of Pacific Biosciences Sequel II system and ultralong reads of Oxford Nanopore. *GigaScience*. 9.
5 328 doi: 10.1093/gigascience/giaa123.
6
7
8
9 329 Mapleson D et al. 2017. KAT: a K-mer analysis toolkit to quality control NGS datasets and genome
10 330 assemblies. *Bioinformatics*. 33:574–576. doi: 10.1093/bioinformatics/btw663.
11
12
13 331 Marçais G, Kingsford C. 2011. A fast, lock-free approach for efficient parallel counting of occurrences of
14 332 k-mers. *Bioinformatics*. 27:764–770. doi: 10.1093/bioinformatics/btr011.
15
16
17 333 Mascher M et al. 2017. A chromosome conformation capture ordered sequence of the barley genome.
18 334 *Nature*. 544:427–433. doi: 10.1038/nature22043.
19
20
21 335 Pfeifer M et al. 2013. The perennial ryegrass GenomeZipper: targeted use of genome resources for
22 336 comparative grass genomics. *Plant Physiol*. 161:571–582. doi: 10.1104/pp.112.207282.
23
24
25 337 Russo A, Potente, Giacomo, Mayionade, Baptiste. 2021. HMW DNA extraction from diverse plants
26 338 species for PacBio and Nanopore sequencing. doi: 10.17504/protocols.io.5t7g6rn.
27
28
29 339 Seppey M, Manni M, Zdobnov EM. 2019. BUSCO: Assessing Genome Assembly and Annotation
30 340 Completeness. In: *Gene Prediction: Methods and Protocols*. Kollmar, M, editor. *Methods in Molecular*
31 341 *Biology* Springer: New York, NY pp. 227–245. doi: 10.1007/978-1-4939-9173-0_14.
32
33
34
35 342 Shafin K et al. 2020. Nanopore sequencing and the Shasta toolkit enable efficient de novo assembly of
36 343 eleven human genomes. *Nature Biotechnology*. 38:1044–1053. doi: 10.1038/s41587-020-0503-6.
37
38
39 344 Smit A.F.A., Hubley R., Green P. RepeatMasker. <http://repeatmasker.org>.
40
41
42 345 Tang H et al. 2015. ALLMAPS: robust scaffold ordering based on multiple maps. *Genome Biology*. 16:3.
43 346 doi: 10.1186/s13059-014-0573-1.
44
45
46 347 Van Bel M et al. 2012. Dissecting Plant Genomes with the PLAZA Comparative Genomics Platform.
47 348 *Plant Physiology*. 158:590–600. doi: 10.1104/pp.111.189514.
48
49
50 349 Walker BJ et al. 2014. Pilon: An Integrated Tool for Comprehensive Microbial Variant Detection and
51 350 Genome Assembly Improvement. *PLOS ONE*. 9:e112963. doi: 10.1371/journal.pone.0112963.
52
53
54 351 Wang Y et al. 2012. MCScanX: a toolkit for detection and evolutionary analysis of gene synteny and
55 352 collinearity. *Nucleic Acids Res*. 40:e49. doi: 10.1093/nar/gkr1293.
56
57
58
59
60

- 353 Wilkins PW, Humphreys MO. 2003. Progress in breeding perennial forage grasses for temperate
354 agriculture. *The Journal of Agricultural Science*. 140:129–150. doi: 10.1017/S0021859603003058.
- 355 Wu Z-Q, Ge S. 2012. The phylogeny of the BEP clade in grasses revisited: Evidence from the whole-
356 genome sequences of chloroplasts. *Molecular Phylogenetics and Evolution*. 62:573–578. doi:
357 10.1016/j.ympev.2011.10.019.
- 358 Zhou Q et al. 2020. Haplotype-resolved genome analyses of a heterozygous diploid potato. *Nature*
359 *Genetics*. 52:1018–1023. doi: 10.1038/s41588-020-0699-x.

360

